



De l'occurrence à l'isotopie. Les co-occurrences en lexicométrie

Damon Mayaffre

► To cite this version:

Damon Mayaffre. De l'occurrence à l'isotopie. Les co-occurrences en lexicométrie. Syntaxe et Sémantique, 2008, 9, pp.53-72. hal-00551114

HAL Id: hal-00551114

<https://hal.science/hal-00551114>

Submitted on 3 Jan 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Syntaxe
&
Sémantique

9

*Textes, documents numériques, corpus.
Pour une science des textes instrumentée*

2008

Presses universitaires de Caen



De l'occurrence à l'isotopie.

Les co-occurrences en lexicométrie

Damon MAYAFFRE

Laboratoire BCL, Nice Sophia-Antipolis
CNRS – MSH de Nice
mayaffre@unice.fr

Résumé : Cette contribution cherche à définir la co-occurrence comme une forme minimale du *co(n)texte*; *co(n)texte* nécessaire à l'élaboration du sens et condition de l'interprétation. Constater que *a* et *b* sont co-occurents n'est rien d'autre que contextualiser l'un par l'autre. Ainsi, nous montrons que l'enjeu du calcul des co-occurrences est essentiel pour la lexicométrie : c'est lui qui, par la contextualisation des formes qu'il propose, fait de la lexicométrie une pratique lexicologique (*versus* lexicographique) et herméneutique (*versus* logico-formelle).

Summary : *This paper attempts to define co-occurrence as a minimal form of co(n)text, namely the context that is a condition for the elaboration of meaning as well as for interpretation. Realizing that a and b are co-occurents merely amounts to contextualising the former through the latter. We shall thus show that working out co-occurrences is crucial in lexicometry. The importance of contextualised forms makes lexicometry a truly lexicological – as opposed to lexicographic-and hermeneutic – as opposed to logico-formal- practice.*

Introduction¹

La lexicométrie – plus généralement l'Analyse des données textuelles – ne s'est jamais présentée comme une pratique *lexicographique*. S'y trouvant souvent ramenée, elle s'est pourtant toujours comprise comme une pratique *lexicologique*.

Si la lexicométrie prit un temps le nom de *lexicologie quantitative* (dans le titre de [Dugast, 1979] par exemple), c'est en effet que ses intentions étaient arrêtées. La statistique et l'informatique furent dès le départ mobilisées non pas seulement pour construire des listes de mots, des index de formes, des dictionnaires de fréquences qui, fussent-ils constitués *en corpus* (et non *en langue*), renvoyaient à une démarche lexicographique. Elles furent mobilisées, fondamentalement, pour nourrir

1. Cet article est directement issu de notre réflexion aux 5^e Journées de Linguistique de Corpus (Lorient, 13 au 15 septembre 2007), <http://web.univ-ubs.fr/corpus/>.

une science du vocabulaire. La co(n)textualisation des unités traitées en devenait alors un élément majeur puisqu'il ne saurait y avoir de réflexion linguistique accomplie sur le vocabulaire (et par-delà sur les textes) sans l'étude étroite des *usages*, c'est-à-dire des *co(n)textes d'utilisation*.

Cette prétention lexicologique dont la contextualisation² est la condition se caractérise par deux types de comportement de l'analyste et deux modes de fonctionnalités classiques des logiciels d'Analyse des données textuelles : le *retour au texte* et le développement d'une *statistique contextualisante*, syntagmatique ou co-occurrence.

(i) Dans les logiciels reconnus sur le marché scientifique tels Hyperbase, Lexico, Astartex ou Weblex, le retour direct au texte plein ou le retour indirect, partiel mais organisé au texte par l'intermédiaire de concordances est systématique à toutes les étapes du traitement. Ainsi pourra-t-on par simple clic accéder au texte intégral *via* une liste de spécificités ou une constellation de mots disposés sur une analyse factorielle des correspondances, et, par là, naviguer dans le corpus pour le lire de façon naturelle. De la même manière, on le sait, il sera possible de convoquer en un instant toutes les phrases ou tous les paragraphes contenant telle ou telle forme puis organiser ces *concordances* (voir récemment [Pincemin *et al.*, 2006]) et en permettre une lecture aisée. Les occurrences – seules et désincarnées – sont donc les entrées utiles et nécessaires pour le traitement lexicométrique mais le retour au (co)-(n)texte est posé comme la condition de l'interprétation. Pour en revenir à des considérations terminologiques, notons ici que les termes en usage (*lexicométrie*, *textométrie*, *logométrie*) peuvent oblitérer, par la force de leur suffixe, cette dimension seulement documentaire, *qualitativiste* ou contextualisante qui représente pourtant une des deux faces du traitement lexicométrique. Grâce à l'hypertextualité, les logiciels organisent des parcours de lecture : ceux-ci sont certes originaux au sens où ils font appel à des ressources informatiques et à une navigation inaccessible aux pratiques manuelles-oculaires, mais l'acte de lecture est ici lui-même traditionnel au sens d'une confrontation entre le texte dans sa chaîne contextuelle naturelle et l'analyste. La lexicométrie est le bras armé de l'*herméneutique numérique* qui voit aujourd'hui le jour (par exemple Viprey 2005a et 2005b ou Mayaffre 2002a et 2006) :

2. Ont été notés ci-dessus « co(n)textualisation » et « co(n)texte ». Renonçons désormais, sauf exception, à alourdir le texte des parenthèses, étant entendu qu'on traitera toujours du *contexte linguistique* aussi appelé *co-texte*. Plus précisément, c'est le *co-texte immédiat* qui sera le plus souvent considéré, bien que la contextualisation linguistique d'une forme ne s'arrête pas à la phrase ou au paragraphe et s'élargit au texte et au corpus (cf. *infra*).

organiser le retour au texte pour en permettre la lecture et favoriser l'acte final interprétatif est une de ses tâches fondamentales.

(ii) Mais retourner trop vite à une lecture traditionnelle du texte, c'est renoncer trop tôt aux traitements quantitatifs dont la lexicométrie postule la pertinence pour lire, comprendre, interpréter les grands corpus textuels. Pour cette raison, le traitement statistique occurrence-tiel d'essence lexicographique – les occurrences nucléaires, décontextualisées qui une fois comptées, triées, indexées sont censées *faire référence* et renvoyer à des ontologies – se prolonge par des traitements statistiques ou mathématiques contextualisants de type co-occurrence-tiel et d'essence lexicologique. En amont des travaux en cours de (Mellet et Barthélemy 2007) ou (Luong, Longrée et Mellet 2008) sur la *topologie textuelle* dont les prétentions de modélisation de la textualité sont plus importantes, c'est de ce traitement co-occurrence-tiel que cette contribution veut traiter, après les travaux pionniers, en France, de (Demonet *et al.*, 1975), (Tournier, 1980), (Lafon, 1984) et ceux plus récents de (Viprey 1997 et 2005a et 2005b), (Heiden 1998 et 2004), (Véronis 2003 et 2004), (Martinez 2003) ou (Brunet 2006, 2007 et 2008)³.

1. De la co-occurrence. Précisions terminologiques

Le terme « co-occurrence » serait bien établi en ADT, si certains auteurs, venant d'autres horizons, ne s'appliquaient à en brouiller le sens : la co-occurrence est la co-présence ou *présence simultanée* de deux unités linguistiques (deux mots par exemple ou deux codes grammaticaux) au sein d'un même contexte linguistique (le paragraphe ou la phrase par exemple, ou encore une fenêtre arbitraire⁴). Cette co-occurrence peut être grossièrement constatée, puis vainement exprimée, en fréquence absolue. Mais plus pertinemment, la lexicométrie la constate et l'exprime grâce à des coefficients statistiques à même de mesurer le degré de significativité des co-présences ou attractions trouvées. Nombre de modèles et de coefficients ont été à ce jour proposés : (Lafon, 1984), (Church et Hanks, 1990), (Dunning 1993), (Fung et

3. Les contributions sur la co-occurrence se heurtent au problème bibliographique : depuis (Firth, 1957) et (Harris, 1957), les articles traitant directement ou indirectement de la co-occurrence sont trop nombreux pour être synthétisés. Nous nous excusons de la multiplication des références qui accompagneront le propos.

4. On le devine, la définition de la taille et de la nature de cette fenêtre contextuelle déterminera l'analyse. Favorisera-t-on un contexte dit naturel telle la phrase ou un contexte arbitraire (fenêtre coulissante de x mots) ? Favorisera-t-on un contexte large comme le paragraphe voire la page, ou un contexte plus étroit comme le syntagme ?

McKeown 1997), (Manning et Schütze 1999), (Véronis, 2003 et 2004), (Wu et Zhou 2003), etc. Et nous rappelons en note à la suite de [Brunet, 2007 et 2008] le mode de calcul hypergéométrique d'influence saint-clousienne implémenté dans Hyperbase qui sera utilisé ici⁵.

Au moins deux termes complémentaires ou concurrents (la *collocation* et la *corrélation*) permettent de préciser les choses pour souligner la dimension générique de la co-occurrence.

Co-occurrence et collocation – Le terme « collocation » apparaît parfois comme synonyme de co-occurrence, particulièrement dans la littérature anglo-saxonne (voir [Williams 1999] et [Daille et Williams 2001] pour diverses définitions). Les collocations pointent pourtant le plus souvent des co-occurents d'un certain type, ceux qui entretiennent des relations syntaxiques (ou parfois distributionnelles). C'est en ce sens que (Hausmann 1979) ou (Mel'čuk *et al.* 1984), pour ne citer que les exemples les plus célèbres, les utilisent. C'est ainsi que les définissent [Béjoint et Thoirion 1992, p. 517] :

Les collocations sont des associations privilégiées de quelques mots (ou termes) reliés par une structure syntaxique et dont les affinités syntagmatiques se concrétisent par une certaine récurrence en discours.

Par là, si la co-occurrence prétend constater statistiquement les usages individuels – des associations libres relevant du choix du locuteur –, la collocation nous renvoie linguistiquement déjà du côté des *contraintes* du système (ou plutôt *des* systèmes idiomatiques) : indéniablement, certaines co-occurrences apparaissent comme des faits de langue, évidemment dans des lexies composées (*chemin de fer*) qui sont hors du champ mais aussi dans les syntagmes semi-figés (*semi-fixed combinations*) objets favoris des études collocatives (*pluie battante, salaire de misère, gravement malade*, etc.)⁶. De fait, la recherche des collocations – qui apparaît donc comme une sous-espèce spécialisée de celle des co-occurrences – est le plus souvent tendue vers la mise à jour, à finalité linguistique, des expressions idiomatiques, des unités phraséologiques, des phrasèmes ou semi-phrasèmes, des locu-

5. Soit s = nombre de phrases ou de paragraphes, f = fréquence du mot-pôle dans le texte, g = fréquence du mot co-occurent dans le texte et k = co-occurrence observée. Alors : $\text{Prob}(x = k) = (f! (s + g)! g! (f + s)! / (k! (f - k)! (g - k)! (s + k)! (f + g + s)!)$

6. Les collocations sont donc pour la plupart des auteurs cette zone floue, ce stade intermédiaire entre associations libres et combinaisons figées. Elles se caractérisent notamment par le fait que chaque élément garde une certaine autonomie et son sens individuel, mais que la combinaison n'est plus, pour le locuteur, tout à fait libre, et produit une plus-value sémantique. Entre langue et discours, c'est cette notion de *semi-figement* ou de *liberté contrainte* qui intrigue les analystes (cf. *infra*).

tions, etc.⁷. Nous retrouvons ainsi la notion de collocation particulièrement présente dans les travaux de traductologie (automatique) afin de déterminer des formules propres à une langue qui exigent un effort de traduction particulier loin du mot à mot. Dans ce cadre, la recherche de collocations aboutit en général aux antipodes de notre propos c'est-à-dire à des travaux lexicographiques, avec l'établissement de nomenclatures ou de dictionnaires censés consigner des locutions, leur sens définitif et éventuellement leur traduction dans d'autres langues : ainsi pourra-t-on consulter pour l'anglais, l'allemand ou le français, les dictionnaires de (Benson *et al.*, 1997), (Ilgenfritz *et al.*, 1989), (Mel'Cuk *et al.*, 1984) ou, plus strictement pour les mots composés, le DELAC de Silberztein et (Gross, 1996).

De manière très générale constatons donc pour conclure que certains praticiens de la co-occurrence ressentent le besoin, au-delà des dépendances syntaxiques strictes collocatives, de *contraindre* ou d'*informer* la relation de co-occurrence statistique. Cela peut passer, par exemple, par une volonté de distinguer le contexte gauche et le contexte droit du mot-pôle ; la valeur de la co-occurrence de *a* et *b* serait différente selon que *a* se trouve à gauche ou à droite de *b*. Cela peut passer encore par la volonté de définir strictement, au sein de la fenêtre contextuelle, l'empan existant entre les deux co-occurents : plus l'empan est réduit plus la co-occurrence serait signifiante.

Toutes ces tentatives, et bien d'autres, pour intéressantes qu'elles soient (elles représentent en effet toutes un raffinement de la méthode), recèlent un danger, celui de pervertir le sens même de la démarche co-occurentielle *stricto sensu* qui consiste à laisser le soin à la statistique de repérer sans a priori les associations privilégiées : que celles-ci soient contraintes ou non, que celles-ci révèlent une dépendance syntaxique ou non, que celles-ci marquent une relation distributionnelle ou non, tout ceci relève pour nous de l'explication du phénomène co-occurentiel mais non de la recherche et du questionnement. Effectivement, dans les mailles du filet co-occurentiel s'accrocheront des lexies composées, des syntagmes figés, des expressions semi-figées mais d'autres associations lexicales plus libres et plus inattendues seront aussi remontées à la surface *sans que l'on veuille a priori ni les exclure ni les privilégier*⁸.

7. Les termes de « phraséologie » ou de « phrasème », au sens ici de Mel'Cuk, nous invitent à préciser que le champ d'investigation des collocations est le plus souvent la phrase (voire le syntagme) lorsque celui de la co-occurrence gagne à être transphrastique (ou a-phrastique) et aime à s'étendre au paragraphe.

8. La collocation est donc synonyme de *semi-figement* ou de combinaisons *semi-figées*. On soulignera, pour conclure, la prudence, dans la bouche même de ceux qui l'utilisent, de la notion de *semi-figement*. Et on rappellera le passage – prudent lui aussi –

Co-occurrence et corrélation – Le terme « corrélation » pourrait présenter l'avantage de suggérer l'approche statistique (cf. les *indices de corrélation* en vigueur dans tout modèle statistique). Là où la co-occurrence de deux termes (si elle est exprimée naïvement en valeur absolue) peut être marginale et fortuite, leur corrélation semble souligner l'intensité de la relation ; la corrélation serait ainsi une co-occurrence *significative* d'un point de vue statistique.

Pourtant, dans la littérature (par exemple Bourion 2001, p. 1), la corrélation pointe une autre réalité. Comme la collocation, la corrélation stigmatise des co-occurents d'un certain type, ceux qui entretiennent une relation sémantique. Deux corrélats seraient deux co-occurents qui ont une relation de sens.

Outre le fait que de manière intriquée – donc problématique – un corrélat (sémantique) peut être un collocat (syntaxique) et vice versa, la notion de corrélation (comme celle de collocation) présente un danger de confusion épistémologique entre les différents plans de l'analyse : nous passons en effet imperceptiblement du *constat* de co-présence statistique (la co-occurrence) à la *signification* linguistique de cette relation (le corrélat). Précisons bien : notre propos immédiat consistera à montrer que la co-occurrence porte en elle un potentiel important pour la sémantique de corpus, la science et l'interprétation du texte : c'est en passant de l'occurrence à la co-occurrence que la lexicométrie accède à la lexicologie et que l'ADT entre dans la sémantique interprétative. Mais nous tenons à distinguer ce qui relève de la description formelle ou matérielle d'un phénomène – description difficilement contestable devant la finesse des outils statistiques et informatiques – et le sens toujours négociable à donner à celui-ci : parler de *corrélats sémantiques* à propos de *co-occurrences statistiques*, c'est conclure ce que nous voulons ici postuler.

2. La co-occurrence : sa dimension herméneutique et ses enjeux pour l'ADT

La lexicologie est l'étude des vocabulaires en usage (Eluerd 2000) : la contextualisation des vocables en est la clef. La linguistique est l'étude des textes définis comme le seul objet empirique du linguiste ([Adam 1999] et [Rastier 2001]) : la contextualisation des unités ou grandeurs textuelles, au sein de parcours de lecture contrôlés, en est ici encore la clef.

Si elle veut servir le vocabulaire et le texte, la lexicométrie doit donc proposer des outils pour traiter du contexte ; le contexte étant

du Cours : « Mais il faut reconnaître que dans le domaine du syntagme, il n'y a pas de limite tranchée entre le fait de langue, marque de l'usage collectif, et le fait de la parole, qui dépend de la liberté individuelle » (Saussure [éd.] 1995, p. 173).

défini autant comme un environnement matériel bien circonscrit (une fenêtre) que comme un moment où le sens prend forme ou un lieu virtuel, éventuellement discontinu, où la textualité prend corps.

Au palier supérieur – Au palier supérieur, le contexte est non seulement *tout le texte* – selon l'expression de François Rastier⁹ –, mais encore le corpus textuel dans son ensemble, macro-objet qui informe ses composants : c'est en effet, pour finir, au sein du corpus que les mots, les phrases, les textes prennent sens pour l'analyste ; c'est au sein du corpus que s'explicitent et s'organisent des stratégies de lecture interprétatives¹⁰.

À vrai dire, cette affirmation – la place centrale du corpus – qui paraît faire aujourd'hui l'unanimité en linguistique (voir récemment, même pour la phonologie, le point de vue de [Laks, 2008]) est LE postulat originel des pratiques lexicométriques. Fondamentalement, la lexicométrie et l'ADT se sont présentées comme une alternative complémentaire à la linguistique introspective en donnant aux linguistes les moyens qui leur manquaient de mettre à l'épreuve leur intuition et d'observer de manière critique les grands corpus empiriques afin de repérer les régularités linguistiques, les caractéristiques ou anomalies du langage réel : il y a là une posture essentielle sur laquelle nous ne pouvons revenir. Techniquement surtout, c'est bien le corpus dans son ensemble qui constitue la norme statistique sans laquelle aucun décompte ne ferait sens. Le dispositif de traitement tout entier repose, on le sait, sur l'idée de *norme quantitative endogène* au corpus¹¹. La fréquence locale d'un mot dans une partie du corpus est mise en rapport avec la

-
9. Par exemple : « Le sens d'une unité est déterminé par son contexte. Le contexte c'est tout le texte : la micro-sémantique dépend donc de la macro-sémantique. » (Rastier 2008).
 10. Nous ne reprendrons pas ici la réflexion sur le rôle et la place du corpus en linguistique. Nous renvoyons le lecteur à la revue dédiée à la question : *Corpus*. Résumons pour les sciences du texte : les corpus sont des objets construits – donc critiques et problématisés – qui informent leurs composants. Nous avons souligné que cette information était d'autant plus performante que les corpus étaient *réflexifs* (Mayaffre 2002b et 2006) c'est-à-dire susceptibles d'internaliser leurs ressources interprétatives : en miroir, les textes du corpus doivent s'éclairer mutuellement ; se *réfléchir* les uns les autres ; chacun d'entre eux constituant le co-texte immédiat de tous, et l'ensemble, l'intertexte de chacun.
 11. Impossible ici aussi de revenir sur ce postulat fondamental : il n'existe pas de fréquence en langue mais seulement en corpus. La nécessité statistique de travailler sur un corpus clos et étalon s'appuie, précisément, sur une conscience lexicologique. Le vocabulaire ne peut être abordé qu'en usage c'est-à-dire *en corpus* : notre objet est le vocabulaire du corpus non le lexique du dictionnaire. *Endogène* est sans doute le mot-clef de la linguistique de corpus, ainsi peut-on envisager une statistique endogène, une stylistique endogène (Viprey 1997), une lexicologie endogène, une sémantique endogène, etc.

fréquence totale dans le corpus. C'est seulement par la mise en contrastes des parties du corpus et sur le postulat que l'ensemble du corpus représente une norme ou un étalon cohérent (une « moyenne » pourrait-on dire rapidement) que repose le traitement quantitatif. Si certains linguistes hésitent encore à admettre que le sens est différentiel, la statistique, elle, ne peut fonctionner sans cette affirmation : on ne jugera la fréquence du mot « autorité » dans le discours de Sarkozy durant la campagne électorale 2007 comme importante, signifiante, porteuse de sens, qu'au regard de la norme que propose par exemple le corpus des discours de tous les candidats à l'élection présidentielle.

Au palier inférieur – Le contexte, linguistique comme statistique, est donc au palier supérieur le corpus. Au palier inférieur, nous voulons poser que *le contexte minimal d'un terme est la co-occurrence*. Nous considérerons en effet, en corpus, que la forme minimale du contexte d'un terme, nécessaire à sa compréhension-interprétation, n'est pas le syntagme ou la phrase mais la co-occurrence ; ou dit autrement encore nous définirons ici la co-occurrence comme la forme minimale du contexte qui présente l'avantage de se trouver accessible de manière systématique, étant entendu que nous saurions considérer, même avec un concordancier, un par un, tous les mots dans toutes leurs chaînes. Pour Saussure [(éd.) 1995, 150 sq.], chaque occurrence est un hapax dont la valeur diffère avec le contexte, l'intonation, etc.¹² : comment faire alors pour réaliser une synthèse des usages ou même seulement relever la moindre régularité ? Pour Guiraud (Guiraud, 1960, p. 19), le sens d'un mot « se définit finalement par la somme de ses emplois » : mais comment faire pour *sommer* des emplois linguistiques (surtout lorsqu'il s'agit d'hapax) ? À mi-chemin, le traitement des co-occurrences entend considérer tous les mots en leurs (con) textes, et extraire de manière systématique de ceux-ci les formes significativement associées à ceux-là ; ou encore, considérer tous les paragraphes du corpus et y repérer systématiquement les associations linguistiques récurrentes jugées comme significatives et, pour cette raison, actrices principales de la textualité.

12. On connaît en effet le passage du Cours récemment rappelé par (Rastier et Valette, à paraître) : « Lorsque, dans une conférence, on entend répéter à plusieurs reprises le mot *Messieurs* ! on a le sentiment qu'il s'agit chaque fois de la même expression, et pourtant les variations de débit et l'intonation la présentent, dans les divers passages, avec des différences phoniques très appréciables [...] ; en outre, ce sentiment de l'identité persiste bien qu'au point de vue sémantique non plus il n'y ait pas identité absolue d'un *Messieurs* ! à l'autre... [...]. Chaque fois que j'emploie le mot *Messieurs*, j'en renouvelle la matière ; c'est un nouvel acte phonique et un nouvel acte psychologique. Le lien entre les deux emplois du même mot ne repose ni sur l'identité matérielle, ni sur l'exacte similitude des sens. » (Saussure 1995, p. 150-152).

Concrètement en tout cas, constater que *a* et *b* sont co-occurents n'est rien d'autre, pour nous, que contextualiser minimalement l'un par l'autre ; déterminer l'ensemble *b, c, d, e*, etc. des co-occurents de *a*, c'est définir l'ensemble des contextes minimaux (mais pertinents) de *a* au sein du corpus.

Enfin, signalons que la définition de la co-occurrence comme forme particulière du contexte aboutit à un glissement de la vision habituelle du contexte ne serait-ce que parce qu'il ne s'agit plus ici de chaîne, de fenêtre, de suite continue mais d'*associations* ; associations discontinues le plus souvent, parfois transphrastiques. Ici la définition matérielle du contexte – le contexte c'est avant tout un co-texte entendu comme un environnement textuel immédiat et contigu – se trouve équilibrée par une dimension herméneutique – le contexte, *c'est ce qui fait sens* ; ce qui sémantise un terme en autorisant l'interprétation. La co-occurrence ainsi définie serait alors une forme indiquée pour incarner le *passage* – passage minimal évidemment – dans sa double dimension que lui attribue (Rastier 2007). La co-occurrence s'exerce d'abord sur le plan du signifiant du texte : pour reprendre les termes de Rastier, la co-occurrence en tant que *passage* est donc un « *extrait de l'expression* » (*ibid.*, p. 30). La chose est essentielle pour nous car le traitement lexicométrique s'appuie toujours sur l'expression, sur des formes bien définies, et ne saurait compter autre chose que des réalités linguistiques matérielles ou implémentables du texte¹³. Pour forcer le trait donc, les co-occurrences sont d'abord, techniquement, des extraits (des mots) extraits (par un algorithme) d'extraits (des paragraphes). Sur le plan du signifié du texte maintenant, la co-occurrence est aussi « un *fragment du contenu* » (*ibid.*), un fragment qui pointe sémantiquement « vers ses contextes gauche et droit, proches et lointains » (*ibid.*). Ce pointage sémantique est médiatisé par la statistique qui est là pour déceler les associations qui font sens *puisqu'on ne saurait les attribuer au hasard* : la co-occurrence statistique, en tant que telle pour l'analyste, produit du sens en proposant une grille de lecture ou un prisme interprétatif du corpus.

Plus simplement, concluons que la co-occurrence, forme minimale du contexte, forme particulière du passage, présente l'avantage d'être

13. Du mot graphique préconisé par St. Cloud aux enchaînements syntaxiques, en passant par les lemmes ou les codes grammaticaux, ces unités sont certes aujourd'hui de plus en plus construites. Nous ne reviendrons pas ici sur la question épineuse des unités de traitement pertinentes en lexicométrie. Posons simplement *a minima* que, dans une posture inductive matérialiste, le traitement informatique lexicométrique doit s'appuyer sur l'expression ou le signifiant avant de traiter du contenu et du signifié.

à la fois une donnée objective de l'énoncé et, par les vertus de la statistique, un objet signifiant du texte.

3. Des co-occurrences de « patrie » chez Sarkozy

Essayons d'illustrer la réflexion théorique par une étude de cas. L'analyse porte sur le discours de la campagne électorale française de 2007. Deux corpus seront traités, le premier est composé des principaux discours des principaux candidats du premier tour à l'élection présidentielle : Laguiller, Buffet, Royal, Bayrou, Sarkozy et Le Pen. Le second contient l'exhaustivité des discours de meeting – 34 exactement – de Sarkozy du 1^{er} janvier 2007 à la veille de son élection ; l'ensemble de ces discours, représentant plus d'un million de mots, est disponible sur le site *Discours 2007* de Jean Véronis (<http://sites.univ-provence.fr/veronis/Discours2007/>).

L'objectif est donc de stabiliser, *via* la recherche des co-occurrences, des contextes minimaux d'utilisation d'un terme afin de déceler des isotopies sémantiques et nourrir l'interprétation.

3.1. Constat occurrenceiel

La démonstration part d'un premier constat lexicométrique, d'ordre seulement occurrenceiel : la distribution du mot « patrie » dans le corpus du premier tour (cf. fig. 1).

Ce constat est suggestif : plus on est à droite de l'échiquier politique plus on emploie « patrie », plus on est à gauche moins on l'utilise. Dans cette grille de lecture politique du corpus, seul Bayrou fait exception. Néanmoins, on ne saurait ici aller plus loin dans l'interprétation textuelle et l'inférence socio-linguistique sauf à supposer trop vite un sens *déjà-là* à « patrie », là où l'on sait que le mot possède un des signifiés les plus problématiques de la langue politique française.

3.2. Le sens de « patrie » : rappel historique

Il en va en effet ainsi de quelques termes dans le discours politique. « Patrie », comme « peuple » par exemple, est dans la bouche d'un homme politique contemporain, à l'image des éléments des rêves, sémantiquement *surdéterminé*. Sans doute peut-on parler de polysémie linguistico-politique en dépit du *Petit Robert* qui accorde une seule entrée et un seul sens au mot (fig. 1).

Tout au contraire, dans le bien nommé *Dictionnaire des usages socio-politiques* (Guilhaumou et Monnier, 2006) montrent que dès l'origine moderne révolutionnaire, le mot, en discours, se charge d'acceptions différentes. On en distingue habituellement trois :

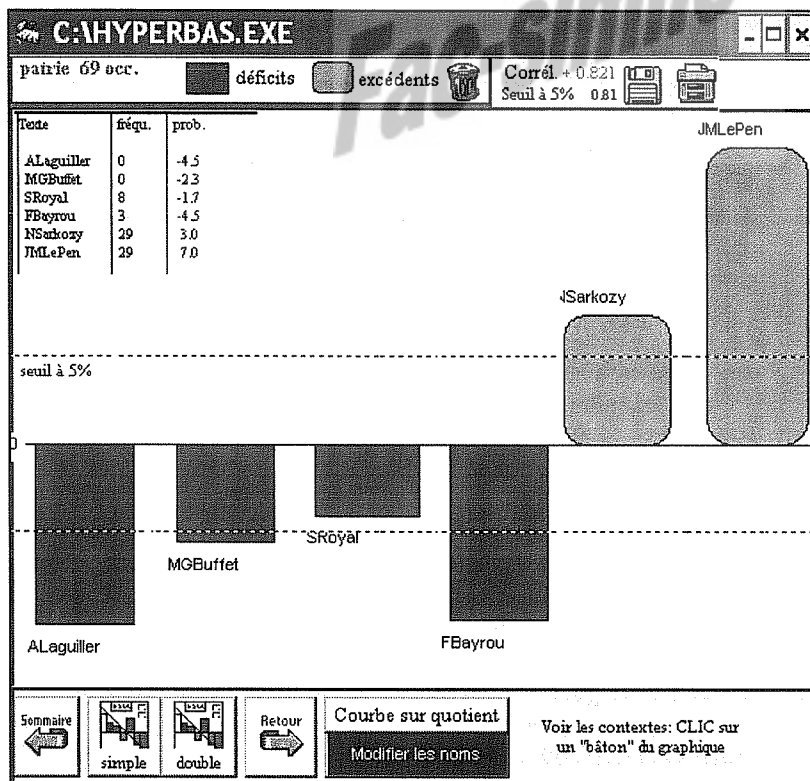


Fig. 1 : Distribution de « patrie » durant la campagne électorale^a

- a. Les sous-utilisations et sur-utilisations mesurées par rapport à la norme du corpus sont ici exprimées en écart réduit. Les non-spécialistes de lexicométrie ne s'étonneront pas, par exemple, que la sous-utilisation de « patrie » soit plus marquée chez Laguiller (-4,5) que chez Buffet (-2,3) alors même que l'une et l'autre n'utilisent pas le terme (fréquence = 0). Cela s'explique par le fait que le sous-corpus de Laguiller est plus important que celui de Buffet : l'absence de « patrie » en devient plus significative.

– une acception territoriale : la patrie c'est le territoire, la frontière à défendre contre les armées étrangères à partir de 1792, la terre des pères (voire la race).

– une acception politique : la patrie, en France, c'est la république des sans-culottes *versus* la monarchie ; par là ce sont des valeurs politiques comme la démocratie, la liberté, l'égalité, la vertu. Loin de la terre, la patrie c'est donc l'Idée (l'idée républicaine s'entend).

– une acception sociale enfin, en lien avec l'acception précédente : la patrie ce sont les crève-la-faim, la paysannerie pauvre et le peuple contre l'aristocratie, les riches, les privilégiés de la noblesse. La patrie

des révolutionnaires ce n'est pas seulement l'idée républicaine naissante mais une revendication de justice sociale.

Dans le « allons enfants de la patrie » de la Marseillaise (25 avril 1792), dans *l'événement discursif* (Guilhaumou 2006) majeur que constitue la déclaration par la Législative de « la patrie en danger » (11 juillet 1792), dans le « vive la patrie ! » de l'armée révolutionnaire de Valmy partant au combat (20 septembre 1792), ces trois dimensions sont présentes. À l'inverse, l'épisode de Coblenz représente, sur ces trois points, le symbole de l'anti-patrie.

Cette polysémie, originelle donc, ne fait ensuite que se complexifier au fil du temps. Tout au long du XIX^e et du XX^e siècle, le mot prend une épaisseur historico-sémantique insondable.

On retiendra, comme signe le plus révélateur de la complexité sémantique du mot, le fait qu'il puisse être aussi bien revendiqué par la gauche que par la droite. En diachronie peut-être peut-on s'entendre sur le fait que né à gauche durant la période révolutionnaire, il passe à droite avec la montée des nationalismes au cours du XIX^e : de part et d'autre, la patrie de Robespierre ne signifie pas la même chose que celle de Pétain. Le basculement s'effectuerait en 30 ans avec le boulangisme, l'affaire Dreyfus, la guerre de 14-18 et la révolution de 1917. Mais ce serait là trop simplifier l'histoire chaotique d'un mot secoué par les grands mouvements de l'histoire (Valmy et Coblenz donc, les guerres napoléoniennes et la Restauration, la guerre avec la Prusse et la Commune, les ligues nationalistes et la Grande guerre, le front populaire et Munich, la Résistance et Vichy, la guerre froide et les guerres coloniales, la construction de l'Europe et la mondialisation).

Bref, sur ce substrat, il est facile de comprendre que l'occurrence seule de « patrie » dans le discours de Sarkozy ne signifie rien ; et le recours au *Robert* ou au *Larousse* risquerait de nous égarer définitivement. Seule, ici comme ailleurs, l'approche contextualisante, *en corpus*, peut instruire le débat d'autant que la plupart des sens historiques qu'un dictionnaire peut consigner ont muté dans la France du XXI^e siècle. La lexicologie est une pratique endogène à un corpus ou elle n'est pas.

3.3. Co-occurrences de « patrie » et isotopies du discours de Sarkozy

Nous avons rappelé récemment, après d'autres, quelques approches différentes du traitement des co-occurrences (Mayaffre 2008) : de l'extraction par le calcul des spécificités des co-occurents d'un mot-pôle donné au repérage systématique de toutes les associations privilégiées du corpus, de la co-occurrence simple à la *poly-cooccurrence* (Martinez 2003), la *Q-occurrence* (Massonnie 1986) ou la *cooccurrence généralisée* (Viprey 1997), tout est aujourd'hui possible pour

pressentir les réseaux thématiques et les isotopies d'un texte, aborder la textualité ou la texture, mettre à jour la cohérence ou le maillage d'un texte conçu comme un entrelacement lexical. Et nous renvoyons le lecteur à la thèse de (Viprey 1997) pour la réflexion théorique la plus aboutie sur l'intérêt de la co-occurrence pour la linguistique textuelle.

Comme l'a montré récemment (Brunet 2006, 2007 et 2008), le logiciel Hyperbase s'applique aujourd'hui à offrir plusieurs outils complémentaires pour embrasser le phénomène. Nous utiliserons, à propos de « patrie » chez Sarkozy, seulement l'outil le plus récent implémenté dans le logiciel, la fonction « Associations » qui permet de constituer *des graphes de co-occurents*.

L'idée de graphe n'est pas nouvelle puisque Demonet *et al.* (Demonet *et al.* 1975) l'avaient pressentie ; Heiden (Heiden 2004) en fait un outil majeur du logiciel Weblex susceptible de produire un *lexicogramme* de l'ensemble des associations du corpus, et Véronis (Véronis 2003 et 2004) propose, avec des applications instructives, l'outil Hyperlex à même de *cartographier* les co-occurrences du texte.

Dans cette lignée, et quoique de facture plus modeste, les graphes d'Hyperbase représentent une concrétion technique de l'idée de *réseaux lexicaux*. Pour cela, le logiciel met en forme le mot-pôle, ses principaux co-occurents statistiques, mais encore les co-occurents de ceux-ci (ou co-occurents indirects), proposant ainsi une profondeur d'analyse à trois niveaux. Par là, le traitement statistique et sa mise en forme graphique donnent à voir des faisceaux isotopiques non triviaux qui se caractérisent, comme on le sait, par des phénomènes quantitatifs de récurrence¹⁴ et d'échos sémantiques complexes que le lecteur peut percevoir dans la trame du texte (cf. figure 2).

Une fois expurgé des mots outils et recentré sur les seuls substantifs, le traitement fait donc ressortir que « patrie » a 5 grands co-occurents¹⁵ chez Sarkozy, qui marquent, si l'on veut bien considérer leurs co-occurents respectifs, cinq dimensions du discours :

14. Dans les définitions qu'en donnent Greimas, Rastier, Arrivé ou Kerbrat, un phénomène isotopique est toujours produit par une *récurrence*, une *redondance*, des *reprises*, parfois des *itérations*. Qu'il nous soit donc permis de constater que l'isotopie fait partie de ces nombreux concepts qui impliquent, sans toujours se l'avouer, un traitement quantitatif. Dit plus directement : comment mesurer raisonnablement une *récurrence* significative sans lexicométrie ?
15. En réalité, les co-occurents dépassant le seuil statistique sont bien plus nombreux. Nous les avons réduits aux cinq substantifs majeurs pour ne pas encombrer le graphique. De la même manière, les substantifs ont été privilégiés pour les co-occurents indirects. Ne cachons pas que la mise en graphe, pour des raisons techniques, réclame toujours des sélections problématiques ; sans rien dire des choix sémiotiques mis en œuvre.

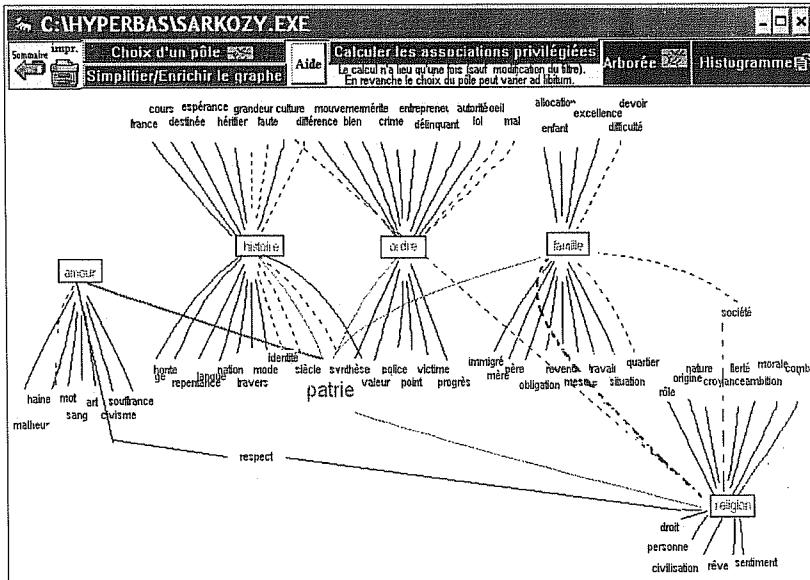


Fig. 2 : Graphe des co-occurents de « patrie » dans le discours de Sarkozy en 2007^a

- a. Encadrés, les co-occurents directs de « patrie ». Non-encadrés, les co-occurents des co-occurents (ou co-occurents indirects de « patrie »). L'épaisseur des traits (de gras à pointillé) indique la force de l'attaction statistique entre deux mots (par simple clic sur le trait le coefficient est indiqué).

a) une dimension pathétique (« patrie » => « amour » -> « sang », « haine », « souffrance », etc.) comme dans cet exemple¹⁶ :

Comment s'étonner qu'en dénigrant l'AMOUR de la PATRIE on réveille le nationalisme qui est la HAINE des autres ? Comment s'étonner que la mode exécrable de la repentance, en voulant faire expier aux Français les fautes supposées des générations passées, ressuscite des HAINES ancestrales que l'on croyait à tout jamais appartenir à l'histoire et rouvre des BLESSURES que le temps avait à peine commencé à fermer ? (Sarkozy, 18 mars 2007, meeting du Zénith à Paris).

Ou, citant Michelet :

Comment être à Rouen et ne pas penser à ce que Michelet disait de Jeanne : « Souvenons-nous toujours, Français, que la PATRIE chez nous est née du

16. Notons que les passages que nous allons donner sont caricaturaux. En effet, ils concentrent le mot-pôle, ses co-occurents directs et les co-occurents indirects. Pourtant, le graphique, lui, n'indique pas que le mot-pôle a un lien nécessaire avec les co-occurents indirects.

CŒUR d'une femme, de sa tendresse, de ses LARMES, du SANG qu'elle a donné pour nous » ? (Sarkozy, 24 avril 2007, meeting de Rouen).

De fait, le discours électoral de Sarkozy, au-delà du cas particulier de « patrie », est un discours qui joue autant sur l'émotion que sur la raison. Nous avons montré ailleurs (Mayaffre, 2007) qu'il se caractérise notamment par son extrémité lexicale (« haine », « détestation », « barbarie », « rêve », « excision », « voyou », etc.) là où le discours républicain classique est un discours de l'euphémisme lexical (« événement », « mouvement de protestation », « délinquant », etc.).

b) Une dimension historique/patriotique (« patrie » => « histoire » -> « France », « grandeur », « culture », « destinée » etc.) ; comme dans cet exemple :

Je me fais une haute idée de la FRANCE, de ce qu'elle incarne aux yeux du monde, de son intelligence, de sa CULTURE, de sa vocation universelle. J'ai fait mienne son HISTOIRE. Pour moi il n'y a pas une HISTOIRE de France de gauche et une HISTOIRE de France de droite. Il n'y en a qu'une parce qu'il n'y a qu'une seule France. J'assume tout, je prends tout en partage et j'en suis fier. Je suis fier d'être un enfant de la PATRIE de Saint Louis, de Voltaire, de Victor Hugo, de Jaurès, de Blum, du Général de Gaulle, de Schuman, de Monnet. (Sarkozy, 11 février, meeting de Versailles).

c) Une dimension politique/autoritaire (« patrie » => « ordre » -> « autorité », « délinquant », « crime », « police », etc.) ; comme ici :

À bas l'AUTORITE ! Cela voulait dire : l'obéissance de l'enfant à ses parents, c'est fini ! Démodé ! La supériorité du maître sur l'élève, c'est fini ! Ringard ! La soumission à la LOI, c'est fini ! Dépassé ! Le pouvoir de POLICE, c'est fini ! Enfin ! Le respect de l'État et de ceux qui le représentent, c'est fini ! L'amour de la PATRIE, la fidélité à la France, à son drapeau, la gratitude vis-à-vis de ceux qui se sont battus pour elle, c'est fini ! La morale, c'est fini ! (Sarkozy, 23 février, meeting de Perpignan).

d) Une dimension sociétale/familiale (« patrie » => « famille » -> « père », « mère », mais aussi « travail », etc.) :

Oui, à force de tout détester, la FAMILLE, la PATRIE, la religion, la société, le TRAVAIL, la politesse, la courtoisie, l'ordre, la morale. À force de tout détester, on finit par se détester soi-même. (Sarkozy, 5 avril 2007, meeting de Lyon).

Peu évoquée par les commentateurs, cette dimension familiale et éducative est l'un des aspects majeurs du discours de Sarkozy (que le calcul des spécificités par exemple révèle bien) et il est seulement étonnant de toucher à ce thème *via* les co-occurrences de « patrie ».

e) Enfin une dimension religieuse/spirituelle (« patrie » => « religion » -> « croyance », « rêve », etc.) que les discours du président Sarkozy sur Dieu et la foi le 20 décembre 2007 à Rome ou en janvier 2008 à Ryad révéleront plus encore.

La contextualisation de « patrie » par ses co-occurents directs puis indirects permet donc sinon de donner un sens définitif au mot¹⁷ en tout cas de retrouver des isotopies endogènes au corpus, assez loin des sentiers sur lesquels nous mènerait la définition dictionnaire. Par exemple, Sarkozy ne développe pas frontalement, en ce début du XXI^e siècle, la question de la patrie-territoire. « Patrie » est un élément isotopique complexe du discours, s'inscrivant dans plusieurs faisceaux distincts, le plus souvent mobilisé (et mobilisateur), dans le cadre de développements idéologiques généraux sur les valeurs comme l'ordre, l'autorité, la morale, le travail. Concrètement, concluons que statistiquement associé aux mots « amour » (co-occurent direct) puis au mot « sang » (co-occurent indirect) par exemple, ou encore associés au mot « famille », lui-même associé au mot « travail », « patrie » se contextualise ; ces associations co-occurentielles constituent en elles-mêmes des contextes minimaux du corpus, ou des passages élémentaires, qui permettent de *lire* le texte, au sens plein, c'est-à-dire de produire du sens et de l'interprétation.

Conclusion

En définissant la co-occurrence comme forme minimale du contexte, cette contribution a essayé d'illustrer, de plusieurs points de vue, une idée unique : le passage d'une approche occurrentielle des corpus textuels à une approche co-occurentielle ne représente pas seulement, pour l'ADT, un saut quantitatif mais une rupture qualitative.

Si la recherche d'occurrences renvoie à une démarche lexicographique ou encore à une linguistique logico-grammaticale dans laquelle il y aurait seulement des entités nucléaires indexées dans un dictionnaire et des règles de composition consultables dans une grammaire, la pratique des co-occurrences est d'essence contextualisante et ouvre d'autres perspectives. Elle réfléchit d'une part une science du vocabulaire *en usage* (l'usage minimal de *a* serait qu'il est, en corpus, statistiquement associé à *b*) et témoigne d'autre part de l'organisation textuelle

17. On a compris au terme de cet article, que, précisément, nous ne croyons ni à la singularité du sens ni à une acception définitive des mots. Les mots sont tous potentiellement polysémiques (notamment en diachronie) et l'étude de leur usage témoigne de cette pluralité sémantique.

en rendant compte, dans une perspective sémantique et herméneutique, de son maillage lexical.

En une formule proposée récemment par [Valette, 2008], la co-occurrence serait alors un élément essentiel d'une « *lexicologie textuelle* », c'est-à-dire d'une science du vocabulaire qui opérerait de manière endogène au texte pour contribuer à en rendre compte.

En cela, le traitement des co-occurrences, souvent abandonné au TALN dans le cadre de la traductologie automatique ou de la désambiguïsation sémantique, est un enjeu majeur pour l'ADT et la lexicométrie qui derrière ses techniques statistiques espère être partie prenante des arts, des sciences et de l'interprétation des textes en proposant des parcours de lecture post-impressionnistes à la fois originaux et contrôlés.

Références bibliographiques

- ADAM J.-M. (1999), *Linguistique textuelle. Des genres de discours aux textes*, Paris, Nathan.
- BEJOIN H. et THOIRON Ph. (1992), « Macrostructure et microstructure dans un dictionnaire de collocations en langue de spécialité », *Terminologie et traduction*, 2-3, p. 513-522.
- BENSON M., BENSON E. et ILSON R. (1997), *The BBI Dictionary of English Word Combinations*, Amsterdam – Philadelphie, John Benjamins Publishing Company.
- BOURION E. (2001), *L'aide à l'interprétation des textes électroniques*, thèse de doctorat, Université de Nancy II (http://www.revue-texto.net/Corpus/Publications/Bourion/Bourion_Aide.html).
- BRUNET E. (2006), « Navigation dans les rafales », in *JADT'06*, J.-M. Viprey (éd.), Besançon, Presses universitaires de Franche-Comté, vol. I, p. 15-29.
- BRUNET E. (2007), « Séquences et fréquences. Mises en œuvre dans Hyperbase », *Lexicométrica*, numéro thématique : *Topographie et topologie textuelles* (<http://www.cavi.univparis3.fr/lexicométrica/numspeciaux/special9/brunet.pdf>).
- BRUNET E. (sous presse), « Les séquences (suite) », in *Actes des JADT 2008* Heiden S. (éd.).
- CHURCH K. W. et HANKS P. (1990), « Word Association Norms, Mutual Information, And Lexicography », *Computational Linguistics*, vol. 16 (1), p. 177-210.

- DAILLE B. et WILLIAMS G. (dir.) (2001), *Collocation: computational extraction, analysis and exploitation*, Workshop at the 39th Annual Meeting and 10th Conference of the European Chapter of ACL, Nantes, Institut de recherche en Informatique.
- DEMONET M. et al. (1975), *Des tracts en mai 1968*, Paris, Colin.
- DUGAST D. (1979), *Vocabulaire et discours: fragments de lexicologie quantitative: essai de lexicométrie organisationnelle*, Genève – Paris, Slatkine – Champion.
- DUNNING T. (1993), « Accurate Methods for the Statistics of surprise and Coincidence », *Computational Linguistics*, vol 19 (1), p. 61-74.
- ELUERD R. (2000), *La lexicologie*, Paris, Puf.
- FUNG P. et MCKEOWN K. (1997), « A technical word and term translation aid using noisy parallel corpora across language groups », *Machine Translation*, 12 (1-2), p. 53-87.
- FIRTH J. (1957), « A Synopsis of Linguistic Theory 1930-1955 », *Studies in Linguistic Analysis*, p. 1-32.
- GROSS G. (1996), *Les expressions figées en français: noms composés et autres locutions*, Paris, Ophrys.
- GUILHAUMOU J. et MONNIER R. (éd.) (2006), *Dictionnaire des usages socio-politiques (1770-1815)*, Paris, Honoré Champion, fasc. 8: « Patrie, patriotisme. Notions pratiques ».
- GUIRAUD P. (1960), *Problèmes et méthodes de la statistique linguistique*, Paris, Larousse.
- HAUSSMANN F. J. (1979), « Un dictionnaire des collocations est-il possible? », *Travaux de linguistique et de littérature*, XVII (I), p. 187-195.
- HARRIS Z. S. (1957), « Co-occurrence and transformation in linguistic structure », *Language*, 33, p. 283-340.
- HEIDEN S. (2004), « Interface hypertextuelle à un espace de cooccurrences: implémentation dans Weblex », *JADT 2004. Le poids des mots*, G. Purnelle et al. (éd.), Louvain, Presses universitaires de Louvain, p. 577-588.
- HEIDEN S. et LAFON P. (1998), « Cooccurrences. La CFDT de 1973 à 1992 », in *Des mots en liberté, Mélanges Maurice Tournier*, Paris, ENS Éditions, t. 1, p. 65-83.
- ILGENFRITZ, P. et al. (1989), *Langenscheidts Kontextwörterbuch Französisch-Deutsch. Ein neues Wörterbuch zum Schreiben, Lernen, Formulieren*, Berlin – Munich, Langenscheidt.
- LAFON P. (1984), *Dépouillements et Statistiques en Lexicométrie*, Genève – Paris, Slatkine – Champion.

- LAKS B. (sous presse), « Pour une phonologie de corpus », *Journal of French Language Studies*.
- LEBART L. et SALEM, A. (1994), *Statistique textuelle*, Paris, Dunod.
- LONGRÉE D., LUONG X. et MELLET S. (sous presse), « Les motifs : un outil pour la caractérisation topologique des textes », *Actes des JADT 2008*, S. Heiden (éd.).
- MANNING C.D. et SCHÜTZE H. (1999), *Foundations of Statistical Natural Language Processing*, The MIT Press, Cambridge (Mass.).
- MARTINEZ W. (2003), *Contribution à une méthodologie de l'analyse des cooccurrences lexicales multiples dans les corpus textuels*, thèse de doctorat, université de la Sorbonne nouvelle – Paris 3, sous la direction d'A. Salem (dactyl.).
- MASSONIE J.-P. (1986), *Pratique de l'analyse des correspondances*, Besançon, Annales Littéraires de l'Université de Franche-Comté.
- MAYAFFRE D. (2002a), « L'Herméneutique numérique », *L'Astrolabe. Recherche littéraire et Informatique* (<http://www.arts.uottawa.ca/astrolabe/auteurs.htm>).
- MAYAFFRE D. (2002b), « Les corpus réflexifs : entre architextualité et hypertextualité », *Corpus*, 1, p. 51-69.
- MAYAFFRE D. (2006), « Philologie et/ou herméneutique numérique : nouveaux concepts pour de nouvelles pratiques », in *Corpus en Lettres et Sciences sociales. Des documents numériques à l'interprétation*, F. Rastier et M. Ballabriga (éd.), Toulouse, Presses universitaires de Toulouse, p. 15-26. (<http://www.revue-texto.net/1996-2007/Parutions/Livres-E/Albi-2006/Mayaff>).
- MAYAFFRE D. (2007), « Vocabulaire et discours électoral de Sarkozy : entre modernité et pétainisme », *La Pensée*, 352, p. 65-80.
- MAYAFFRE D. (sous presse), « Quand “travail”, “famille”, “patrie” co-occurrent dans le discours de Nicolas Sarkozy. Étude de cas et réflexion théorique sur la co-occurrence », *Actes des JADT 2008*, Heiden (éd.).
- MEL'CUK I. et al. (1984), *Dictionnaire Explicatif et Combinatoire du français contemporain I*, Montréal, Presses de l'Université de Montréal.
- MELLET S. et BARTHÉLEMY J.-P. (2007), « La topologie textuelle : légitimation d'une notion émergente », *Lexicométrie*, numéro thématique « Topographie et topologie textuelle » (<http://www.cavi.univ-paris3.fr/lexicometrica/numspeciaux/special9/mellet.pdf>).
- PINCEMIN B. et al. (2006), « Concordanciers : thème et variations », *JADT'06*, Jean-Marie Viprey (éd.), Besançon, Presses universitaires de Franche-Comté, vol. II, p. 773-784.

- RASTIER F. (2001), *Arts et sciences du texte*, Paris, PUF.
- RASTIER F. (2007), « Passages », *Corpus*, 6, p. 25-54.
- RASTIER F. (sous presse), « Entretien sur les théories du signe et du sens », à paraître en traduction anglaise in F. Stjernfeld et P. Bundgaard (éd.), *Theories of Signs and Meaning*, Automatic Press.
(http://www.revue-texto.net/docannexe/file/1735/bundgaard_rastier.pdf)
- RASTIER F. et VALETTE M. (à paraître), « De la polysémie à la néo-sémie », *Langue française*.
- SALEM A. (1993), *Méthodes de la statistique textuelle*, thèse d'État, thèse de doctorat, université de la Sorbonne nouvelle – Paris 3 (dactyl.).
- SAUSSURE F. de (1995), *Cours de linguistique générale*, Paris, Payot.
- TOURNIER M. (1980), « D'où viennent les fréquences de vocabulaire ? La lexicométrie et ses modèles », *Mots*, 1, p. 189-209.
- VALETTE M. (sous presse), « À quoi servent les lexiques sémantiques généralistes ? Discussion et propositions », *Cahiers du Cental*.
- VÉRONIS J. (2003), « Cartographie lexicale pour la recherche d'information », in *Actes de TALN 2003*, Batz-sur-mer, ATALA IRIN
(http://www.atala.org/doc/actes_taln/AC_0100.pdf).
- VÉRONIS J. (2004), « Hyperlex : lexical cartography for information retrieval », *Computer, Speech and Language*, 18 (3), p. 223-252.
- VIPREY J.-M. (1997), *Dynamique du vocabulaire des Fleurs du mal*, Paris, Honoré Champion.
- VIPREY J.-M. (2005a), « Philologie numérique et herméneutique intégrative », in *Sciences du texte et analyse de discours*, J.-M. Adam et U. Heidmann (éd.), Genève, Slatkine, p. 51-68.
- VIPREY J.-M. (2005-b), « Corpus et sémantique discursive : éléments de méthode pour la lecture des corpus », in *Sémantique et corpus*, A. Condamines (dir.), Paris, Lavoisier, p. 245-276.
- VIPREY J.-M. (2006), « Structure non séquentielle des textes », *Langages*, 163, p. 71-85.
- WILLIAMS G. (1999), *Les réseaux collocationnels dans la construction et l'exploitation d'un corpus dans le cadre d'une communauté de discours scientifique*, thèse de doctorat, université de Nantes (dactyl.).
- WU H. et ZHOU M. (2003), « Synonymous collocation extraction using translation information », in *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, E. Hinrichs et D. Roth (éd.), Sapporo, National Institute of Informatics, p. 120-127.